

# Year 4 — Practical Numerical Analysis

Based on lectures by Dr Kathryn Gillow

Notes taken by James Arthur

Michaelmas 2022

These notes are not endorsed by the lecturers, and I have modified them (often significantly) after lectures. They are nowhere near accurate representations of what was actually lectured, and in particular, all errors are almost surely mine (especially the typos!).

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Root Finding (and optimisation)	2
1.1.1	Univariate Root Finding	2
1.1.2	Regular Falsi	3
1.1.3	Newton Raphson	3
1.1.4	Polynomial Root Finding	4
1.1.5	Multivariate Root Finding	4
<b>2</b>	<b>IVPs</b>	<b>6</b>
2.1	Euler Methods	6
2.2	Trapezium Rule / Crank Nicolson Scheme	7
2.3	Generalisation – $\theta$ -method	7
2.4	Euler via Taylor Series	7
2.5	Truncation Error	7
<b>3</b>	<b>Runge-Kutta Schemes</b>	<b>9</b>
3.1	Derivation of Explicit Runge-Kutta Scheme	9
3.2	Stability	10

# 1 Introduction

Problem Sheet by Tuesday at 12 noon. Thursday on the problem sheets from week 2. TA is Robert McDonald. We will study,

- Rootfinding
- ODEs
  - Euler Schemes
  - Runga Kutta
  - Linear multistep
- Parabolic PDEs - Heat Equation

## 1.1 Root Finding (and optimisation)

The idea is simple. Find some  $x$  such that  $f(x) = 0$ . We are being vague on purpose. The stating is simple, the solving isn't simple. The cubic formula isn't so simple. Then for the quintic, then it's only has closed form when its Galois group is solvable. These closed forms, may not exist, be unstable, require evaluation of special function. Hence we need numerical methods.

The relation to optimisation is quite nice. The max or min occurs at a turning point. That is,

$$f_i := \frac{\partial g}{\partial x_i} = 0$$

and then we have a root finding problem.

### 1.1.1 Univariate Root Finding

In the 1D case, given  $f : [a, b] \rightarrow \mathbb{R}$  then find  $c \in [a, b]$  such that  $f(c) = 0$ . The bisection algorithm uses the intermediate value theorem (See Analysis / Topology). Thus to find a root of  $f(c)$  we find  $a$  and  $b$  such that  $f(a)$  and  $f(b)$  have opposite signs. Then let  $c = (a + b)/2$ . Then compute  $f(c)$ , check and repeat.

We need to decide when to terminate the algorithm, the normal are,

- $|b - a| < \text{tol}$ , or,
- $|f(c)| < \text{tol}$ .

For convergence, since  $c \in [a, b]$  the maximum error is  $b - a$ . Since the step halves, the error at step  $n$ ,

$$|c_n - c| \leq \frac{b - a}{2^n}.$$

Thus to achieve an accuracy of tolerance requires,

$$n \geq \frac{\log(b - a) - \log(\text{tol})}{\log(2)}$$

steps of bisection algorithm.

Here is some pseudocode for this algorithm,

---

**Algorithm 1** Bisection Method

---

**Input:**  $a, b, f(x), \text{tol}$ **Output:**  $c$ 

```

while  $|b - a| > \text{tol}$  do
   $c \leftarrow (a + b)/2$ 
  if  $f(a)f(c) < 0$  then
     $b \leftarrow c$ 
  else
     $a \leftarrow c$ 
  end if
end while

```

---

**1.1.2 Regular Falsi**

This algorithm finds a clever guess. We are going to approximate  $f$  on  $[a, b]$ ,

$$p_1(x) = \frac{x-a}{b-a}f(b) + \frac{b-x}{b-a}f(a)$$

Then  $p_1(x)$  has a root at,

$$c = \frac{af(b) - bf(a)}{f(b) - f(a)}$$

and then we use this in the bisection algorithm.

This fails for functions like  $f(x) = x^{10} - 1/10$ . There is a fix,

$$c = \frac{af(b)/2 - bf(a)}{f(b)/2 - f(a)}$$

This is known as the illinois algorithm.

**1.1.3 Newton Raphson**

Suppose we approximate  $f(x)$  by  $g(x)$  which is the truncated Taylor series of  $f$  about the point  $x_k$ ,

$$g(x) = f(x_k) + (x - x_k)f'(x_k)$$

Then  $g(x)$  has a root at,

$$c = x_k - \frac{f(x_k)}{f'(x_k)}$$

Thus the newton raphson method requires an intial guess and then iterates,

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

until convergence is achieved. This is quicker, quadratic, if we have a guess sufficiently close.

We can variate this to secant method. The iterate is,

$$x_{k+1} = x_k - f(x_k) \frac{x_k - x_{k+1}}{f(x_k) - f(x_{k-1})}$$

This avoids the need for  $f'(x_k)$ . We also have damped newtons,

$$x_{k+1} = x_k - \sigma_k \frac{f(x_k)}{f'(x_k)}$$

**Algorithm 2** Newton-Raphson Method**Input:**  $x_0, f(x), f'(x), \text{tol}$ **Output:**  $x^*$  $k \leftarrow 0$ **while**  $|f(x_k)| > \text{tol}$  **do** $x_{k+1} \leftarrow x_k - f(x_k)/f'(x_k)$  $k \leftarrow k + 1$ **end while** $x^* \leftarrow x_k$ 

for some  $\sigma_k \in (0, 1]$ . A third is Halley's method,

$$x_{k+1} = x_k - \frac{2f(x_k)f'(x_k)}{2(f'(x_k))^2 - f(x_k)f''(x_k)}$$

This works if the roots  $x_*$  satisfies  $f'(x_*) \neq 0$  and is then Newton's method applied to the function  $f(x)/\sqrt{|f'(x)|}$ .

**1.1.4 Polynomial Root Finding**

Suppose we want to find the root of a Polynomial,

$$p(z) = \sum_{i=0}^n a_i z^i$$

where  $a_n = 1$ . Methods based on bisection will only find real roots, one at a time. Methods based on Newton's will find roots one at a time. With a complex initial guess such methods will find complex roots.

In order to find all the roots at once, we can define the companion matrix  $C \in \mathbb{C}^{n \times n}$  as,

$$C = \begin{pmatrix} 0 & & & -a_0 \\ 1 & \ddots & & \vdots \\ & \ddots & 0 & -a_{n-2} \\ & & 1 & -a_{n-1} \end{pmatrix}$$

Then we can consider  $zI - C$ ,

$$zI - C = \begin{pmatrix} z & & & a_0 \\ -1 & \ddots & & \vdots \\ & \ddots & z & a_{n-2} \\ & & -1 & z + a_{n-1} \end{pmatrix}$$

and then we can see that  $\det(zI - C) = p(z)$ . Hence the roots of  $p(z)$  is equivalent to finding the eigenvalues of matrix  $C$ . This is what matlab `roots` uses this approach. You may think this is circular, but we can find eigenvalues in other ways, like the QR algorithm.

**1.1.5 Multivariate Root Finding**

Now we want to find  $x$  such that  $f(x) = 0$  for  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . Algorithms based on bisections aren't extendable to higher dimensions. Newton based algorithms are easier to extend. Again we approximate  $f(x)$  by the first few terms,

$$f(x + \delta x) \approx f(x) + J(x)\delta(x)$$

where  $J$  is the jacobian. Then we need to find the newton iterate,

$$J(x)\delta x = -f(x)$$

such that,

$$J(x_k)(x_{k+1} - x_k) = -f(x_k)$$

We can also write the iterate as,

$$x_{k+1} = x_k - (J(x))^{-1}f(x_k)$$

thus damped can be written as,

$$x_{k+1} = x_k - \sigma_k(J(x))^{-1}f(x_k)$$

Here's the kicker. We don't really know how many solutions a root finding algorithm has. We need to look more closely. Hence we bring the idea of newton fractals. For example, suppose we have  $f(z) = z^3 - 1$ . This has three roots,  $z = 1, (-1 \pm \sqrt{3}i)/2$ . By writing  $z = x + iy$  and equating real and imaginary parts. We get,

$$\begin{aligned} x^3 - 3xy^2 - 1 &= 0 \\ 3x^2y - y^3 &= 0 \end{aligned}$$

which can be solved using Newtons Method. For each point for  $(x, y) \in [-2, 2]$  we compute a solution and record what root it converges to. A newton fractal is what this converges to.

## 2 IVPs

**Theorem 2.1** (Picard). Suppose that  $f(t, u)$  is a continuous function of  $t$  and  $u$  in a region  $\Omega = [0, T] \times [u_0 - \alpha, u_0 + \alpha]$  of the  $(t, u)$  plane and that there exists an  $L > 0$  such that,

$$|f(t, u) - f(t, v)| \leq L|u - v|, \quad \forall (t, u), (t, v) \in \Omega$$

Suppose also that,

$$MT \leq \alpha$$

where  $M = \max_{\Omega} |f|$ . There is a unique continuously differentiable function  $u(t)$  defined on  $[0, T]$  satisfying,

$$\begin{aligned} \frac{du}{dt} &= f(u, t) \\ u(0) &= u_0 \end{aligned}$$

Suppose we want to solve,

$$\begin{aligned} \frac{du}{dt} &= f(u, t) \quad t > 0 \\ u(0) &= u_0 \end{aligned}$$

In order to solve this numerically over  $[0, T]$  we define a set of time points at which we wish to approximate the solution. We set  $t_n = n\Delta t$  for  $n = 0, 1, \dots, N$  where  $\Delta t = T/N$ . Then we can integrate,

$$u(t_{n+1}) = u(t_n) + \int_{t_n}^{t_{n+1}} f(t, u(t)) dt$$

Using different approximations of the integral lead to different numerical scheme

### 2.1 Euler Methods

Let  $U_n$  be the numerical approximation to  $u$  at  $t_n$ . For explicit (or forward) Euler we let,

$$\int_{t_n}^{t_{n+1}} f(t, u(t)) dt \approx \Delta t f(t_n, u(t_n))$$

This gives,

$$U_{n+1} = U_n + \Delta t f(t_n, U_n)$$

or,

$$\frac{U_{n+1} - U_n}{\Delta t} = f(t_n, U_n),$$

for  $n = 0, 1, \dots, N-1$  and  $U_0 = u_0$ .

For implicit Euler, we use

$$\int_{t_n}^{t_{n+1}} f(t, u(t)) dt \approx \Delta t f(t_{n+1}, u(t_{n+1}))$$

and this gives...

Explicit Euler is very simple, we can just computer the approximations. The implicit euler is more complication, as we get a non-linear equation. This then becomes a root finding method, but as the solutions are an iterated solution to an ODE so the initial guess should be just the previous  $U_n$ .

## 2.2 Trapezium Rule / Crank Nicolson Scheme

Another way to approximate is,

$$\int_{t_n}^{t_{n+1}} f(t, u(t)) dt \approx \frac{\Delta t}{2} (f(t_n, u(t_n)) + f(t_{n+1}, u(t_{n+1})))$$

This then gives a numerical scheme,

$$U_{n+1} = U_n + \frac{\Delta t}{2} (f(t_n, U_n) + f(t_{n+1}, U_{n+1}))$$

## 2.3 Generalisation – $\theta$ -method

Both the implicit and explicit euler methods as well as Crank Nicolson are specific  $\theta$ -methods,

$$\frac{U_{n+1} - U_n}{\Delta t} = \theta f(t_{n+1}, U_{n+1}) + (1 - \theta) f(t_n, U_n)$$

and our special cases are,

- $\theta = 0$ , Explicit Euler
- $\theta = 1$ , Implicit Euler
- $\theta = \frac{1}{2}$ , Crank Nicolson.

For all non-zero  $\theta$ , the method is implicit and a non-linear equation must be solved at each time-step.

## 2.4 Euler via Taylor Series

The explicit and implicit euler scheme can be motivated using Taylor series expansions. Consider expanding  $u(t_{n+1})$  about  $t_n$ . We have,

$$u(t_{n+1}) = u(t_n) + \Delta t u'(t_n)$$

and then rearranging and fiddling,

$$\frac{u(t_{n+1}) - u(t_n)}{\Delta t} + \mathcal{O}(\Delta t) = f(t_n, u(t_n))$$

## 2.5 Truncation Error

We have just seen that the Euler methods can be derived by truncating Taylor series. The truncation for the  $\theta$ -method is defined as,

$$T_n = \frac{u_{n+1} - u_n}{\Delta t} - \theta f(t_{n+1}, u_{n+1}) - (1 - \theta) f(t_n, u_n)$$

The truncation error can be computed by finding Taylor series about some point. For  $\theta = 0$ , the expansions are usually about  $t = t_n$ , while for  $\theta = 1$  we usually go for  $t = t_{n+1}$  and for anything else, it doesn't matter so we expand about  $t_{n+1/2} = \frac{t_n + t_{n+1}}{2} = t_n + \Delta t/2$ . For explicit we get,

$$T_n = u'(t_n) - f(t_n, u(t_n)) + \frac{1}{2} \Delta t u''(\tau_n)$$

Then we remember  $u'(t) = f(t, u(t))$  and so,

$$T_n = \frac{1}{2} \Delta t u''(\tau_n).$$

Then for the theta method, it gets worse. Welcome to Taylor series hell...

$$T_n = \frac{\Delta t}{2}(1 - 2\theta)u''(t_{n+1/2}) + \mathcal{O}(\Delta t^2)$$

We see that the  $\mathcal{O}(\Delta t^2)$  never vanishes and so  $T_n$  will never be zero. Thus,

$$\begin{cases} \mathcal{O}(\Delta t) & \text{for } \theta \neq 1/2 \\ \mathcal{O}(\Delta t^2) & \text{for } \theta = 1/2 \end{cases}.$$

More precisely, we can show that,

$$T_n = \begin{cases} \frac{\Delta t}{2}u''(\tau_n^{(1)}) & \theta = 0 \\ -\frac{\Delta t}{12}u'''(\tau_n^{(2)}) & \theta = 1/2 \\ -\frac{\Delta t}{2}u''(\tau_n^{(2)}) & \theta = 1 \end{cases}$$

The order of a method is defined to be  $p$  where  $p$  is the largest integer such that  $T_n = \mathcal{O}(\Delta t^p)$



### 3 Runge-Kutta Schemes

We continue the study of,

$$\begin{aligned} u'(t) &= f(t, u) \\ u(0) &= u_0 \end{aligned}$$

We say if  $\theta \neq 0$  we had to solve a non-linear equation, but we can avoid this with more complicated methods. We also saw that Crank Nicholson was very good on error. We now go from Crank Nicholson, using Explicit Euler. We call this improved euler,

$$\frac{U_{n+1} - U_n}{\Delta t} = \frac{1}{2}(f(t_{n+1}, U_n) + \Delta t f(t_n, U_n) + f(t_n, U_n)).$$

We can write Runge-Kutta as,

$$\frac{U_{n+1} - U_n}{\Delta t} = \sum_{i=1}^s b_i k_i$$

where,  $k_1 = f(t_n, U_n)$  and,

$$k_i = f\left(t_n + c_i \Delta t, U_n + \Delta t \sum_{j=1}^{i-1} a_{i,j} k_j\right),$$

for  $i = 2, \dots, s$ . Then  $k_i$ 's are known as the stages of the method and the method is often referred to as an  $s$ -stage method. We usually put the coefficients of explicit Runge-Kutta scheme in some Butcher Tableaux. The improved Euler Scheme can be written as,

$$\frac{U_{n+1} - U_n}{\Delta t} = \frac{1}{2}(k_1 + k_2)$$

where  $k_1 = f(t_n, U_n)$  and  $k_2 = f(t_n + \Delta t, U_n + \Delta t k_1)$ . Another is the modified Euler-Scheme,

$$\frac{U_{n+1} - U_n}{\Delta t} = f\left(t_n + \frac{1}{2}\Delta t, U_n + \frac{1}{2}\Delta t f(t_n, U_n)\right)$$

The butcher table is,

**Remark.** We note that the **b** entries should sub to 1.

#### 3.1 Derivation of Explicit Runge-Kutta Scheme

The coefficients are chosen to make the methods as high order as possible, we can do this using Taylor series expansions of the truncation error. The truncation error is,

$$T_n = \frac{u(t_{n+1}) - u(t_n)}{\Delta t} = \sum_{i=1}^s b_i \tilde{k}_i$$

where  $\tilde{k}_1 = f(t_n, u(t_n))$  and so on. Now death by taylor series. (I am not typing this. )

**Theorem 3.1.** The order  $p$  of an explicit  $s$ -stage Runge-Kutta method is bounded by  $p \leq s$ . Further, it is possible to construct Runge-Kutta methods that achieve this maximal order.

### 3.2 Stability

$u'(t) = \lambda u$ . If  $\lambda < 0$ , the exact solution limits to 0. The explicit euler scheme can be solved,

$$U_n = (1 + \lambda\Delta t)^n$$

For  $\lambda < 0$  we require it to limit to zero. Hence,

$$-1 < 1 + \lambda\Delta t < 1$$

and so,

$$\Delta t < \frac{2}{|\lambda|}.$$

We can do a similar to improved euler. We solve,

$$U_n = (1 + \lambda\Delta t + \frac{1}{2}(\lambda\Delta t)^2)$$

and we see again,

$$\Delta t < \frac{2}{|\lambda|}.$$

The interval of absolute stability is the interval of values  $\lambda\Delta t$  such that  $\lim_{n \rightarrow \infty} U_n = 0$ . Thus for both explicit and improved Euler,  $\lambda\Delta t \in (-2, 0)$ .